# Journal of Computational Biology

**Editor-in-Chief:** Mona Singh, PhD
ISSN: 1066-5277 | Online ISSN: 1557-8666 | Published Monthly | Current Volume: 28

**Impact Factor:** *1.479 ⓘ*2020 Journal Citation Reports (Clarivate, 2021)

**CiteScore™:** 2.1

ⓘ

The leading peer-reviewed journal in computational biology and bioinformatics, publishing in-depth statistical, mathematical, and computational analysis of methods, as well as their practical impact.

- **View Aims & Scope**

  **Indexing/Abstracting**

- Featured Content
- About This Publication
- Editorial Board
- Reprints & Permissions
- News Releases
- Sample Content

## Aims & Scope

*Journal of Computational Biology* publishes articles whose primary contributions are the development and application of new methods in computational biology, including algorithmic, statistical, mathematical, machine learning and artificial intelligence contributions. The journal welcomes novel methods that tackle established problems within computational biology; novel methods and frameworks that anticipate new problems and data types arising in computational biology; and novel methods that are inspired from studying natural computation. Methods should be tested on real and/or simulated biological data whenever feasible. Papers whose primary contributions are theoretical are also welcome. Available only online, this is an essential journal for scientists and students who want to keep abreast of developments in bioinformatics and computational biology.

**Research Articles:** Research articles describe new methodology development and application in computational biology. It is recommended that manuscripts should be approximately 3,000 words, excluding tables, figures, legends, abstract, disclosure or references; longer articles can also be

submitted. Research articles should include the following sections, in order: abstract, introduction, methods, results, discussion and references.

**Software articles:**  Short 2-4 page articles describing implementations of new or recently developed computational methods for applications in computational biology.   The approaches underlying the software should have methodologically interesting components. Software articles can be published as companion articles to primary research articles which describe the main methodological contributions.   Software article submissions should be accompanied by a cover letter that concisely states the novel implementation and algorithmic challenges the software tackles.

\*Research and software articles should report unique findings not previously published.

**Tutorials:**  These articles highlight important concepts in computational biology. The journal especially welcomes tutorials on algorithms, data structures, machine learning paradigms, and other computational formalisms that are newly being utilized in computational biology.  Prospective contributors should contact the journal (**jcb@liebertpub.com**) with brief outlines before proceeding.

**Reviews:**  Brief outlines from prospective contributors are welcome, and these will also be solicited on specific subjects.   Articles that benchmark existing approaches are also welcome.

**News/Perspectives/Book Reviews:**  These article types should typically be 2–4 pages long.  Contacting the journal before beginning such a paper is suggested.

**Conference and other special issues:** The Journal of Computational Biology welcomes proposals for special issues related to topics within the scope of the journal.

## *Journal of Computational Biology* coverage includes:

- Algorithms for computational biology
- Mathematical modeling and simulation
- AI / Machine learning
- Statistical formulations
- Software for applied bioinformatics
- Genomics and systems biology
- Evolution and population genomics
- Biomedical applications
- Biocomputing and biology-inspired algorithms

**Specific topics of interest include, but are not limited to:**

- Molecular sequence analysis
- Sequencing and genotyping technologies
- Regulation and epigenomics
- Transcriptomics, including single-cell
- Metagenomics
- Population and statistical genetics
- Evolutionary, compressive and comparative genomics
- Structure and function of non-coding RNAs
- Computational proteomics and proteogenomics
- Protein structure and function
- Biological networks
- Computational systems biology
- Privacy of biomedical data
- Bioimaging

*Journal of Computational Biology* is under the editorial leadership of Editor-in-Chief **Mona Singh, PhD**, Princeton University; and other leading investigators. View the entire **editorial board**.

 **Audience:** Computational biologists, bioinformaticians, data scientists, applied mathematicians, and computer scientists, among others.

## Indexing/Abstracting:

- PubMed/MEDLINE
- PubMed Central
- Current Contents®/Life Sciences
- Biotechnology Citation Index®

- Biochemistry & Biophysics Citation Index®
- Biological Abstracts
- BIOSIS Previews
- Journal Citation Reports/Science Edition
- EMBASE/Excerpta Medica
- Scopus
- Chemical Abstracts
- ProQuest databases
- CAB Abstracts
- Global Health
- MathSciNet
- The DBLP Computer Science Bibliography
- BenchSci

## Society Affiliations

**The Official Journal of:**

**RECOMB**

RECOMB

# Journal of Computational Biology

**A Journal of Computational Molecular Cell Biology**

The Official Journal of

## RECOMB

*Mary Ann Liebert, Inc.* publishers

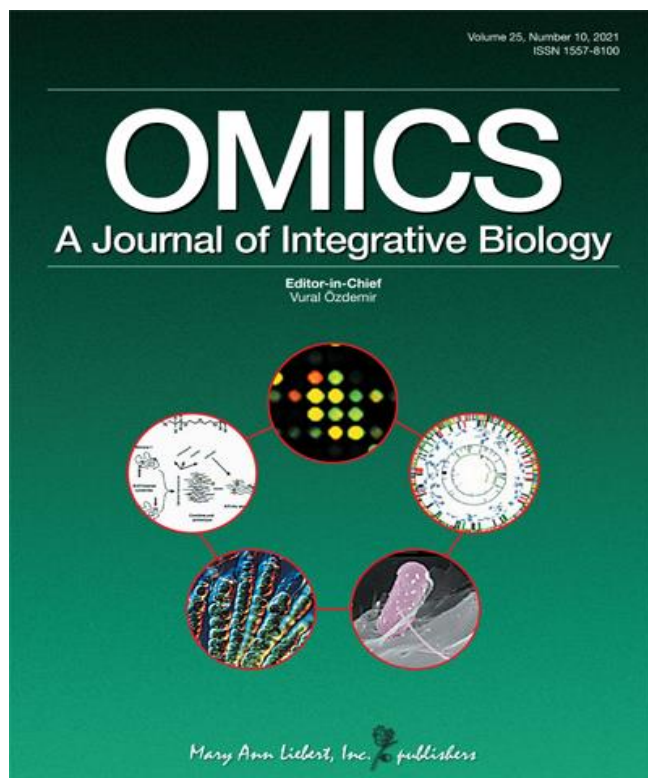2021 February

**Special Issue:**
ICCABS 2019

READ MORE
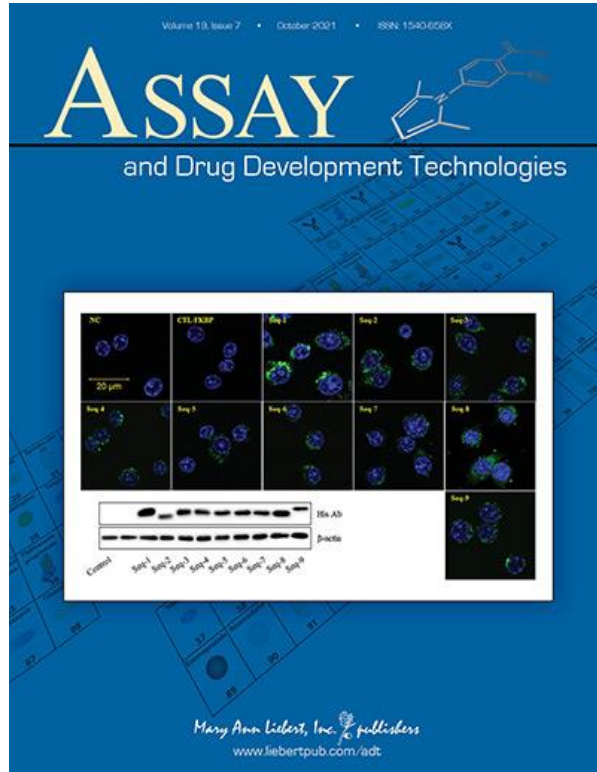
**More Special Issues...**
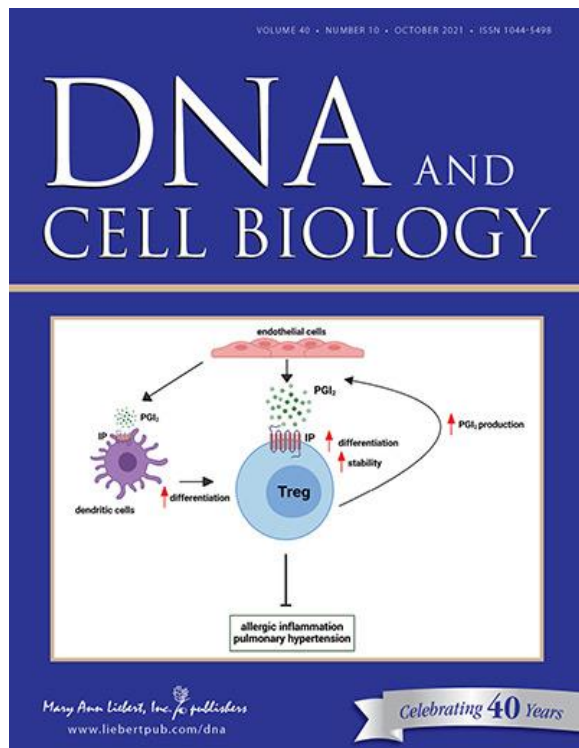
**Recommended Publications**



Genetic Engineering & Biotechnology News

OMICS: A Journal of Integrative Biology



ASSAY and Drug Development Technologies

DNA and Cell Biology



Network and Systems Medicine

**Publications**

- Publications A-Z
- Journal Collections
- Publications by Type
- Recommend a Title

**For Authors**

- Submission Guidelines & Policies
- Author Services
- Fees and Options
- For Reviewers
- Publishing Open Access
- Permissions & Reprints

**Librarians**

- Our Journals
- Account Support
- Archive
- Terms & Conditions
- Resources
- Liebert Link Newsletter

- Contact

**Open Access**

- Open Option
- Open Access Journals
- Publishing Services
- FAQs
- Contact

**Corporate Capabilities**

- 
- Custom Publications
- Interactive Media
- Other Opportunities
- Reprints

- **Advertising**
- **Company**
- **Customer Support**
- **Contact Us**
- **Privacy Policy**

## Bioinformatics Next-generation Whole

Soltani, najma

Department of Computer, Maybod Branch, Islamic Azad University, Maybod, Iran

**Abstract**

It has become progressively obvious that one of the significant obstacles in the genomic age will be the bioinformatics difficulties of cutting edge sequencing. We give an outline of an overall system of bioinformatics examination. For every one of the three phases of (1) arrangement, (2) variation calling, and (3) sifting and comment, we portray the examination required and study the distinctive programming bundles that are utilized. Moreover, we examine conceivable future improvements as information sources develop and feature openings for new bioinformatics devices to be created.

**Keywords**: Computation; Bioinformatics; Sequence; Values.

## 1. Introduction

Without question, the improvement of cutting edge sequencing has changed biomedical examination. Several subsequent age sequencing stages, like Roche/454, Illumina/Solexa, AB/SOLiD, and LIFE/Ion Torrent, have made high-throughput hereditary investigation all the more promptly open to specialists and even clinicians [1]. Not too far off, third era sequencing innovations, like Oxford Nanopore, Genia, NABsys, and GnuBio, will keep on expanding throughput abilities and reduction the expense of sequencing. With each new age of sequencing innovation, there is a remarkable expansion in the surge of information. The genuine difficulties of high throughput sequencing will be bioinformatics. As ever bigger datasets become more moderate, computational examination instead of sequencing will be the rate-restricting component in genomics research. In this paper, we give an outline of the current computational system and choices for genomic examination and give some point of view toward future turns of events and forthcoming requirements.

In this paper, we will examine a portion of the choices in every one of the means and give a worldwide point of view toward the product "pipelines" at present being developed (Figure 1).
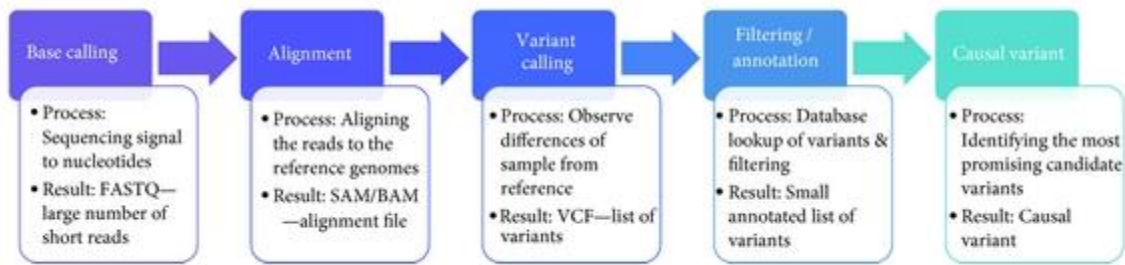
**Figure 1**
Workflow of next-generation sequencing bioinformatics.

## 2. Overview

While distinctive sequencing advancements might utilize diverse starting crude information (e.g., imaging documents or variances in ebb and flow), the possible yields are nucleotide base calls. Short strings of these bases, changing from handfuls to many base sets for each section, are consolidated together, frequently in a type of a FASTQ record. From here, bioinformatics investigation of the succession falls into three general advances: (1) arrangement, (2) variation calling, (3) sifting and comment.

The initial step is arrangement—coordinating with every one of the short peruses to positions on a reference genome (for the motivations behind this paper, the human genome). The subsequent arrangement is put away in a SAM (succession arrangement/guide) or BAM (double arrangement/map) record [2]. The subsequent advance is variation calling—contrasting the adjusted successions and realized groupings to figure out which positions digress from the reference position. The produces a rundown of positions or calls recorded in a VCF (variation call design) document [3]. The third step incorporates both sifting just as explanation. Separating takes the huge number of variations and decreases them to a more modest set. For tumors, this includes contrasting harmful cell genomes with ordinary genomes. For family information, it includes choosing variations that adjust to a particular hereditary legacy design. Explanation includes questioning known data about every variation that is identified. Comment might uncover, for instance, that a variation is a definitely known single nucleotide polymorphism, that a practical impact has effectively been anticipated, that the capacity or action of the quality being referred to is now known, or even that a related sickness has been distinguished.

At last, the ideal outcome from the investigation is few all around commented on variations that can clarify a natural wonder. For instance, for a Mendelian sickness, investigation could recognize the causative variation or quality. For malignancy, investigation might highlight driver changes or targetable qualities. Beginning from base calls and finishing with naturally significant hereditary variations, each progression of investigation might be performed utilizing one of many bits of

programming. This paper examines a few of the bioinformatics choices for every one of these three stages.

## 3. Alignment

Arrangement is the most common way of planning short nucleotide peruses to a reference genome. Since every one of the large numbers of short peruses should be contrasted with the 3 billion potential situations inside the human genome, this computational advance isn't minor. Programming should survey the logical beginning stage of each read inside the reference genome, and the assignment is muddled by the volume of short peruses, special versus non-exceptional planning, and variety in base quality. This progression is consequently computationally serious and tedious [4]. It is additionally a basic advance, as any mistakes in arrangement to the reference genome will be brought through to the remainder of the examination.

The Sequence Alignment/Map (SAM) and Binary Alignment/Map (BAM) designs are the standard record designs for putting away NGS read arrangements [2]. There are different programming programs, some industrially accessible and others unreservedly accessible to established researchers, that can be utilized to perform sequencing read arrangement. Different projects vary in speed and precision. Most arrangement calculations utilize an ordering strategy to all the more quickly limited down potential arrangement areas inside the reference genome with ungapped arrangement, albeit different calculations consider gapped arrangement. Various ways to deal with arrangement include hash tables, separated seeds, as well as bordering seeds. This strategy additionally empowers correlation of varying yield structures (single versus numerous conceivable arrangement yields) [5].

Short peruses produced from NGS may either be single-end peruses or matched end peruses from the example, and may go from handfuls to many base sets [5]; these peruses should be adjusted accurately to their fitting area inside the reference genome. Calculations commonly use a hash-based list (e.g., MAQ, ELAND), BWT-based record (e.g., BWA, Bowtie, SOAP2), genome-based hash (e.g., Novoalign, SOAP), or a separated seed approach (e.g., SHRiMP). A few calculations report the "best" match utilizing heuristic methodologies (e.g., BWA, Bowtie, MAQ), while others take into consideration all conceivable matches (e.g., SOAP3, SHRiMP). Calculations contrast in whether they can deal with both single-end and combined end peruses, or only one sort (e.g., SARUMAN for single-end peruses), and regardless of whether they can perform gapped arrangement (e.g., BWA, Bowtie2) notwithstanding ungapped arrangement (e.g., MAQ, Bowtie). A few calculations center around speed (e.g., BWA, Bowtie), some on affectability (e.g., Novoalign), and a few calculations intend to the two (e.g., Stampy). Table 1 gives a posting of applicable calculations for arrangement of short peruses to the reference genome. While there has been past correlations about these calculations [6], we depict a portion of the more up to date programs, like Bowtie and Bowtie 2, or SOAP/SOAP2/SOAP3, and others beneath.

**Table 1**

### 3.1. Bowtie/Bowtie 2

The Bowtie calculation is both ultrafast and memory productive [7] because of its utilization of a refinement of the FM Index, which itself uses the Burrows-Wheeler change for ultrafast and memory-effective arrangement of peruses to a reference genome. Necktie refines BWT with an original quality-mindful backtracking calculation that grants confounds. In any case, there might be a few tradeoffs among speed and arrangement quality utilizing this calculation [5]. Bowtie2 takes into account investigation of gapped peruses, which might result either from genuine inclusions or cancellations, or from sequencing blunders. The more up to date adaptions use full-text minute records and equipment sped up unique programming calculations to streamline both speed and exactness [8].

### 3.2. BWA/BWA-SW

The BWA approach, in view of BWT, gives effective arrangement of short peruses against the reference genome [9]. This is the most usually utilized methodology for grouping arrangement, and followed the advancement of the original hash-table based arrangement calculation MAQ [10]. BWA refined MAQ by taking into account gapped arrangement of single end peruses, which is significant for longer peruses that might contain indels, and considered sped up. BWA-SW takes into consideration matches without heuristics and arrangement of longer groupings [11].

### 3.3. mrFAST/mrsFAST

As opposed to calculations zeroed in on "interesting" arrangement of areas of the genome and determination of the "best" match, mrFAST [12] and mrsFAST [13] consider quick evaluation of duplicate number variety and task of successions into both one of a kind and the more mind boggling copied locales of the genome [14, 15]. The procedure of these calculations is a seed-and-stretch out approach like BLAST, which uses hash tables to list the reference genome. These calculations can deal with more modest underlying variations (e.g., indels) and bigger primary variations like inclusions, cancellations, reversals, CNVs, and segmental duplications in a reserve neglectful way.

### 3.4. SHRiMP/SHRiMP2

Created to deal with a more noteworthy number of polymorphisms by using a factual model to screen out bogus positive hits, SHRiMP [16] can be used for shading dispersed peruses from AB SOLiD sequencers and can likewise be utilized for normal letter-space peruses. SHRiMP2 [17] empowers direct arrangement for combined peruses and uses different divided seeds, however rather than utilizing recorded peruses like SHRiMP, SHRiMP2 changed to an ordering strategy like Bowtie and BWA.

### 3.5. SOAP/SOAPv2/SOAPv3

Cleanser was created for use in gapped and ungapped arrangement of short peruses utilizing a seed procedure for either single-read or pair-end peruses, and can likewise be applied to little RNA and mRNA label successions [18]. SOAP2 diminished memory utilization and sped up utilizing BWT

for hash-based ordering rather than the seed calculation, and furthermore incorporates SNP recognition [19]. SOAP3 is a GPU (designs handling unit) rendition of the compacted full-text file based SOAP2, which takes into consideration a speed improvement [20].

## 4. Variant Calling

After arrangement of the short peruses to the reference genome, the following stage in the bioinformatics interaction is variation calling. Since the short peruses are as of now adjusted, the example genome can measure up to the reference genome and variations would then be able to be distinguished. These variations might be answerable for illness, or they may essentially be genomic commotion with no practical impact. Variation call design (VCF) is the normalized nonexclusive arrangement for putting away grouping variety including SNPs, indels, bigger underlying variations and comments [3]. The computational difficulties in SNP (variation) calling are because of the issues in recognizing "valid" variations versus arrangement and additionally sequencing mistakes. However the capacity to distinguish SNPs with both high affectability and explicitness is a vital stage in recognizing succession variations related with illness, discovery of uncommon variations, and evaluation of allele frequencies in populaces.

The trouble of variation calling is convoluted by three elements: (1) the presence of indels, which address a significant wellspring of bogus positive SNV distinguishing pieces of proof, particularly if arrangement calculations don't perform gapped arrangements; (2) mistakes from library readiness because of PCR antiquities and variable GC content in the short peruses except if combined end sequencing is used; and (3) variable quality scores, with higher blunder rates commonly found at bases at the finishes of peruses [4]. Along these lines, the pace of bogus positive and bogus negative calls of SNVs and indels is a worry. A nitty gritty audit of SNP-calling calculations and difficulties suggests recalibration of per-base quality scores (e.g., GATK, SOAPsnp), utilization of an arrangement calculation with high affectability (e.g., Novoalign, Stampy), and SNP calling utilizing Bayesian strategies or probability proportion tests and consolidation of linkage disequilibrium to further develop SNP call exactness [21]. We give an outline of a portion of the product bundles for variation calling beneath.

### 4.1. The Genome Analysis ToolKit (GATK)

Created by the Broad Institute, the Genome Analysis ToolKit (GATK) is one of the most famous strategies for variation calling utilizing adjusted peruses. It is planned in a measured manner and depends on the MapReduce practical programming approach [22]. The bundle has been utilized for activities, for example, The Cancer Genome Atlas [23] and the 1000 Genomes Project [24] that have covered investigations of HLA composing, different grouping realignment, quality score recalibration, numerous example SNP genotyping and indel revelation and genotyping [22].

### 4.2. SOAPsnp

Created by the Beijing Genome Institute, SOAPsnp is an open source calculation (http://soap.genomics.org.cn/) that expects admittance to an excellent variation data set utilizing

SOAP arrangement results as an information [18]. It very well may be utilized for agreement calling and SNP identification for the Illumina Genome Analyzer stage and uses the phred-like quality score to work out the probability of every genotype dependent on the arrangement results and sequencing quality scores. Expanding upon the speed of the arrangement calculation Bowtie [7] and utilizing SOAPsnp for SNP calling, an open source distributed computing apparatus called Crossbow [7] was created to perform both arrangement and SNP calling.

### 4.3. VarScan/VarScan2

Created by the Genome Institute at Washington University in St. Louis, VarScan (http://genome.wustl.edu/instruments/disease genomics/) is an open source apparatus for short read variation recognition of SNPs and indels that is viable with different sequencing stages and aligner calculations, for example, Bowtie and Novoalign [25]. It can distinguish variations at 1% recurrence, which can be helpful for pooled tests; VarScan licenses examination of individual examples also. VarScan2 [26] incorporates a few upgrades upon VarScan, for example, the capacity to investigation cancer typical example sets for substantial transformations, LOH (loss of heterozygosity) and CNAs (duplicate number adjustments). This program understands growth and ordinary example Samtools accident or mpileup yield all the while for pairwise correlations of base calling and standardized grouping profundity at each position.

### 4.4. ATLAS 2

Created by the Baylor Genome Center, Atlas 2 can be utilized for variation calling of adjusted information from different NGS stages on a scope of processing stages [27]. Map book 2 can likewise be executed by means of a web asset called Genboree Workbench (http://www.genboree.org/). A couple of other electronic investigation apparatuses are accessible like DNANexus (http://www.dnanexus.com/) and Galaxy [28]. Subtleties of Atlas 2 in contrast with other variation calling calculations, for example, SAMtools mpileup and GATK are remembered for Challis et al. [27], and checked on by Ji [29].

### 5. Insertions and Deletions

While most of examination has zeroed in on sicknesses related with SNPs, indel (addition and cancellation) transformations are a typical polymorphism that can likewise exhibit to natural impacts. Studies have shown that little indels may be profoundly connected with neuropsychiatric illnesses like schizophrenia, chemical imbalance, mental hindrance, and Alzheimer's sickness [30].

Moreover, the presence of certain indels is related with the sickness movement of HBV-initiated hepatocellular carcinoma (HCC) in the Korean populace [31]. Indels are likewise utilized as hereditary markers in normal populaces [32]. With the development of sequencing stages and examination instruments, identification of indels through NGS has become more normal. Nonetheless, precise planning of indels to the reference genome is testing, since it requires approaches that include convoluted gapped arrangement and matched end grouping derivation [9]. Also, the event pace of indels is roughly 8-crease lower than that of SNPs [33]. An ideal blend of

both arrangement and indel-calling calculations is fundamental for recognizing indels with high affectability and explicitness. One audit assessed the presentation of different arrangement devices on microindel recognition, and suggested single-end peruses gapped arrangement planning devices, for example, BWA and Novoalign [34]. Different programming approaches have been created to distinguish indels, including an example development approach (e.g., Pindel) and a Bayesian technique (e.g., Dindel). A nitty gritty audit by Neuman et al. assessed the presentation of a few distinction indel-calling programs within the sight of shifting boundaries (read profundity, read length, indel size, and recurrence). By utilizing both mimicked and genuine information that incorporated the Caenorhabditis elegans genome, they saw that Dindel has the most noteworthy affectability (indels found) at low inclusion, despite the fact that Dindel is just appropriate for Illumina information examination. VarScan and GATK require extra boundary changes, like high inclusion for VarScan, to arrive at their best exhibition. This survey gives data to proper instrument choice and boundary enhancement to help fruitful trial plans and suggests Dindel as an appropriate device for low inclusion tests. Underneath, we review the instruments that have been generally utilized for indel calling.

### 5.1. Pindel

Pindel is a product program which executes an example development way to deal with identify breakpoints of huge cancellations (1–10 kb) and medium estimated inclusions (1 bp–20 bp) from combined end short peruses in NGS information [35]. A new, further developed, variant, Pindel2, has been acquainted which incorporates the capacity with distinguish additions of any size, reversals and couple duplications [35]. Pindel has been utilized for the 1000 Genomes Project (http://www.1000genomes.org/) [36], the Genome of the Netherlands project, and the Cancer Genome Atlas [23].

### 5.2. Dindel

Created by the Welcome Trust Sanger Institute, Dindel is an open-source program that uses a Bayesian methodology for calling little (<50 bp) inclusions and erasures (http://www.sanger.ac.uk/assets/programming/dindel/) [37]. Primarily, this calculation realigns succession peruses planned to an assortment of up-and-comer haplotypes that address elective arrangements to the reference. Dindel has been utilized in the 1000 Genomes Project call sets and can just dissect information from Illumina.

### 5.3. GATK

As portrayed in the variation calling area, the Genome Analysis ToolKit (GATK), which gives an assortment of information examination devices, can likewise permit indel calling dependent on the MapReduce programming approach [22]. Subtleties of GATK in contrast with other indel calling techniques including Dindel (VarScan, SAMtools mpileup) are assessed in Neuman et al. [38].

### 6. Filtering and Annotation

After arrangement and variation calling, a rundown of thousands of possible contrasts between the genome under examination and the reference genome is produced. The subsequent stage is to

figure out which of these variations are probably going to add to the obsessive interaction under investigation. The third step includes a mix of both separating (eliminating variations that fit explicit hereditary models or are absent in ordinary tissue) just as explanation (looking into data about variations and recognizing ones that fit the organic cycle).

Sifting should be possible with a hereditary family or with malignant growth and ordinary examples from a similar person. In the example of malignant growth, a typical technique is eliminating variations that are available in both the disease test and the ordinary example, leaving just substantial variations, which have changed from the germline succession. In the example of a family, separating should be possible dependent on the diverse legacy designs. For instance, if the legacy design is autosomal latent, the variations that are heterozygous in the guardians and homozygous in the kid can be picked. Comparative strategies should be possible with bigger families dependent on the legacy design.

As well as separating, further choice of causal variations can be founded on existing comment or anticipated practical impact. Many apparatuses exist to look at applicable variations by referring to recently known data about their organic capacities and deriving potential impacts dependent on their genomic setting. Also, many instruments have been created to recognize hereditary variations that cause illness pathogenesis or phenotypic change [39]. Uncommon nonsynonymous SNPs will be SNPs that cause amino corrosive replacement (AAS) in the coding locale, which conceivably influence the capacity of the protein coded and could add to sickness.

The development of exome and genomic sequencing is yielding a broad number of human hereditary variations, and various illness related SNVs can be distinguished after arrangement and variation calling. In contrast to gibberish and frameshift changes, which frequently bring about a deficiency of protein work, pinpointing illness causal variations among various SNVs has become one of the significant difficulties because of the absence of hereditary data. For example, 1,300 loci are demonstrated to be related with 200 illnesses by GWASs yet a couple of these loci have been recognized as sickness causing variations [40]. Exome sequencing empowers the ID of more clever hereditary variations than beforehand conceivable, yet it actually requires computational and test ways to deal with foresee whether a variation is pernicious. To this end, a few methodologies have been created to recognize uncommon nonsynonymous SNPs that cause amino corrosive replacement (AAS) in the coding locale. The significant guideline of the protein-succession based techniques to anticipate malice in the coding arrangement depends on relative genomics and useful genomics. Relative sequencing examination expects that amino corrosive deposits that are basic for protein capacity ought to be saved among species and homologous proteins; hence, changes in exceptionally rationed destinations are bound to bring about more malicious impact. Different modalities to anticipate infection causing variations incorporate protein organic chemistry, like amino corrosive charge, the presence of a limiting site, and

construction data of protein. SNVs that are anticipated to adjust protein include (like extremity and hydropathy) and design (restricting capacity and modification of auxiliary/tertiary construction) have a higher likelihood of being pernicious.

Albeit most of exploration has zeroed in on protein-changing variations, noncoding variations establish an enormous part of human hereditary variety. Results got from GWAS show that 88% of characteristic related powerless impact variations are found in noncoding locales, exhibiting the significance of practical explanation of both coding and noncoding variations [41]. Computational apparatuses for protein-arrangement based expectation of injuriousness fall into two classifications: requirement based indicators like MAPP and SIFT, and prepared classifiers like MutationTaster and polyPhen. Notwithstanding protein-arrangement based techniques, one more approach to focus on infection relaxed SNVs is through nucleotide-grouping based forecast in noncoding and coding DNA. This cycle likewise uses relative genomics to anticipate harmfulness, and is utilized by projects like phastCons, GERP, and Gumby. In one itemized survey of sickness causing variation ID, the creators presented the ideas and instruments that permit hereditary comment of both coding and noncoding variations [39]. They additionally looked at the overall utility of nucleotide-and protein-based methodologies utilizing exome information, finding that nucleotide-based limitation scores characterized by Genomic Evolutionary Rate Profiling (GERP) and protein-based malicious effect scores gave by PolyPhen were like two Mendelian illnesses, proposing that nucleotide-based forecast can be pretty much as incredible as protein-based measurements [39]. Underneath, we study devices that are useful recognizing illness causal variations among various up-and-comers.

### 6.1. Sorting Intolerant from Tolerant (SIFT)

Arranging Intolerant From Tolerant (SIFT) (http://sift.jcvi.org/) expectation depends on rationed amino corrosive deposits through various species utilizing near sequencing investigation through PSI-BLAST [42]. This depends on the assumption that amino corrosive buildups that are fundamental for protein capacity ought to be evolutionally saved by normal choice. Consequently, SNPs bringing about AAS on the moderate deposits are bound to be malicious.

### 6.2. PolyPhen

PolyPhen/PolyPhen2 (http://genetics.bwh.harvard.edu/pph2/) calculation predicts the likely effect of AAS on the construction and capacity of human protein dependent on protein arrangement, phylogenetic and primary data [43]. An amino corrosive substitution may happen at a particular site where restricting to different atoms or the development of an auxiliary/tertiary design is disturbed. Hence, PolyPhen decides whether the AAS is found at a site which is clarified as a disulfide bond, a functioning site, a limiting site, or a particular theme, for example, transmembrane area. One more capacity of PolyPhen is to look at the succession and polymorphic districts of homologous proteins in the very family to recognize AASs that are uncommon or never saw in the family. Furthermore, PolyPhen likewise guides of the replacement site to the realized 3-dimensional protein design to evaluate if an AAS can possibly annihilate protein structure by

means of a modification of, for instance, the hydrophobic center of a protein, electrostatic associations, or cooperations with ligands or different atoms.

### 6.3. VariBench

VariBench (http://structure.bmc.lu.se/VariBench/) is the principal benchmark information base that gives testing and preparing instruments to computational variety impact forecast [44]. It contains tentatively approved variety datasets gathered from the writing and applicable data sets. The datasets housed in VariBench empower ID of variations that influence protein resistance, protein solidness, record factor restricting locales, and graft destinations. Also, VariBench maps variation positions to the DNA, RNA, and protein groupings at RefSeq, and to the 3-dimensional protein structures at Protein Data Bank (PDB).

### 6.4. snpEFF

snpEFF is an open source, Java-based program that quickly orders SNP, indel, and MNP variations in genomic groupings as having either high, medium, low or modifier useful impacts [45]. Variation comment depends on genomic area (intron, exon, untranslated district, upstream, downstream, graft site, intergenic locale) and anticipated coding impact (interchangeable/nonsynonymous amino corrosive substitution, acquire/loss of start/stop site, frameshift transformations). The program might track down a few distinct capacities for a solitary variation due to contending expectations dependent on elective records. snpEFF utilizes a VCF info and yield style. Right now snpEFF doesn't uphold underlying variations yet there are plans to fuse such help soon. snpEFF is viable with GATK and Galaxy, which are well known variation calling tool stash. The program right now upholds 260 genome forms and can be utilized with custom genomes and comments.

### 6.5. The SNPeffect Database

The SNPeffect Annotation data set (http://snpeffect.switchlab.org/) utilizes arrangement and design data to foresee the impact of protein-coding SNVs on the primary aggregate of proteins [46]. It is principally centered around illness causing and polymorphic variations in the human proteome. This program looks at variation protein forecasts to wild kind protein data from the UniProtKB information base, which right now contains in excess of 60,000 variation proteins. Variation portrayal is accomplished by coordinating conglomeration, amyloid expectation, chaperone-restricting forecast, and protein solidness investigation data by applying a few calculations to every wild kind and freak protein. The principal calculation, TANGO, identifies areas that are inclined to accumulation and works out a score distinction between the freak and wild sort protein. The WALTZ calculation is applied to foresee amyloid-shaping locales in protein successions utilizing a position-explicit scoring framework to derive amyloid-framing penchant. LIMBO is a calculation that predicts chaperone restricting destinations for the Hsp70 chaperones. In situations where underlying data is accessible, the FoldX calculation is utilized to work out the distinction in free energy between the transformed protein and the wild sort and decide if the transformation settles or undermines the construction. Changes are likewise portrayed as falling

into reactant destinations as per data in the Catalytic Site Atlas or not, and falling into known spaces or not. Subcellular data is anticipated utilizing PSORT.

## 6.6. SeattleSeq

SeattleSeq (http://snp.gs.washington.edu/SeattleSeqAnnotation/) comments on known and novel SNPs with natural capacities, protein positions and amino-corrosive changes, preservation scores, HapMap frequencies, PolyPhen expectations, and clinical affiliations dependent on a coordinated information base. A large portion of the explanation data is gotten from the Genome Variation Server (http://gvs.gs.washington.edu/GVS134/), which incorporates data from dbSNP just as different sources. The calculation acknowledges input records in various configurations including GATK and VCF yield styles. Right now, explanation of indels is restricted.

## 6.7. ANNOVAR

The ANNOVAR programming apparatus (http://www.openbioinformatics.org/annovar/) uses exceptional data to quickly practically comment on hereditary variations called from sequencing information [47]. ANNOVAR chips away at various assorted genomes including hg18, hg19, mouse, worm, fly, and yeast. The comment framework permits the client adaptability in the arrangement of genomic areas that are questioned. Comments can be quality based (clients can choose the quality definition framework; RefSeq, UCSC, ENSEMBL, GENCODE, and so on), locale based (record factor restricting destinations, DNAse I excessive touchiness locales, ENCODEmethylation destinations, segmental duplication locales, DGV destinations, and so forth), channel based (e.g., utilizing just variations revealed in dbSNP, or just variations with MAF > 1%), or in light of any of numerous other client driven functionalities.

## 6.8. The Variant Annotation, Analysis and Search Tool (VAAST)

The Variant Annotation, Analysis and Search Tool (VAAST) recognizes harmed qualities and pernicious variations in close to home genome groupings utilizing a probabilistic hunt technique [48]. The apparatus uses both existing amino corrosive replacement and aggregative ways to deal with variation prioritization and consolidates them into a solitary brought together probability system. This strategy builds the precision with which illness causing variations are recognized. VAAST scores both coding and noncoding, and both uncommon and normal, variations all the while and totals this data to distinguish sickness causing variations.

## 6.9. The Variant Analysis Tool (VAT)

The Variant Analysis Tool, VAT, (http://vat.gersteinlab.org/) practically comments on variations called from individual genomes at the record level and gives synopsis measurements across qualities and people [49]. Tank is a computational system that can be executed through an order line interface, a web application, or a virtual machine in a distributed computing climate. This instrument has been used widely to explain loss-of-work variations acquired as a component of the 1000 Genomes Project [50]. The VAT modules snpMapper, indelMapper and svMapper relate SNPs, indels and SVs to protein-coding qualities while the genericMapper module relates variations to noncoding locales of the genome. Record level investigation permits ID of influenced

isoforms. Tank yields VCF records just as representation summing up the organic effect of explained variations.

## 6.10. VARIANT

Variation (VARIant ANalysis Tool) (http://variant.bioinfo.cipf.es/) gives comment of variations from cutting edge sequencing dependent on a few distinct data sets and archives including dbSNP, 1000 Genomes Project, the GWAS inventory, OMIM, and COSMIC [36]. The gave comments likewise remember data for the administrative or underlying jobs of the variations just as the specific pressing factors on the influenced genomic locales. In contrast to other such devices, VARIANT uses a far off information base and works by interfacing with this data set through effective RESTful Web Services. At present VARIANT backings all human, mouse and rodent qualities. Examining variations created by exome sequencing of families in which uncommon Mendelian illnesses are isolated can be a tedious cycle.

## 6.11. VAR-MD

VAR-MD is a product device to break down variations got from exome or entire genome sequencing in human families with Mendelian legacy [51]. This calculation yields a positioned rundown of potential illness causing variations dependent on anticipated pathogenicity, Mendelian legacy models, genotype quality, and populace variation recurrence information. This apparatus is remarkable in that it utilizes family-based comment of grouping information to improve transformation distinguishing proof. VAR-MD is a Unix-based instrument and is executed in Python. Free elements of the program are typically run successively. To work with equal investigation of numerous informational collections, VAR-MD uses Galaxy for dispersed asset the executives.

The different variation comment devices vary in the sorts of variations they measure. All calculations cycle SNPs and indels, however a couple, like ANNOVAR and VAT, can deal with SVs. These devices likewise vary in the registering climate wherein they are executed. Some depend on order line activity while others work utilizing online interfaces or virtual machines in the cloud. A few apparatuses use neighborhood data sets while other utilize around date far off data sets. These different apparatuses likewise vary in the genomic districts that they target. For instance, SNPeffect centers around the proteome while different apparatuses center around the more subtle, yet at the same time practically pertinent areas. From the extensive rundown of potential variations, through sifting and explanation, a more modest rundown of most likely causal variations is produced.

## 7. Conclusion

While the current apparatuses in every one of the three phases of the bioinformatics examination are satisfactory, more information will empower further huge upgrades. New innovation and calculations may fundamentally move the field unforeseeablely, however a few future upgrades

are unsurprising as (1) sequencing peruses expansion long, (2) additional genomes are finished, and (3) comment data sets are better populated.

In the first place, as sequencing innovation expands the base pair read length, arrangement will turn out to be more precise. More limited peruses match with a more prominent number of genome destinations. As peruses fill long, they can be planned all the more exactly with less alternatives and accordingly less wiggle room. This is particularly obvious in areas with low intricacy or a high number of rehashes, traditionally undeniably challenging locales to plan. Longer peruses will make arrangement a simpler issue.

Second, the course of variation calling will profit from bigger data sets of finished genomes. A variation is gotten from correlation with the reference genome, and our arrangement of reference genomes keeps on developing. This will empower variation calling dependent on ethnic foundation, or in view of populaces of genomes rather than a solitary reference genome or a little arrangement of reference genomes.

Third, while separating shows up far-fetched to change fundamentally, explanation and useful expectation will be improved by more information and more-populated data sets. For separating, since the hereditary models and expulsion of typical variations from cancer variations depend just on the hereditary qualities and the examples under investigation, extra data from the data sets won't change these angles a lot. Conversely, the viability of explanation is straightforwardly identified with what is available in known data sets. Various elements of information, for example, useful, pathway, biochemical, or hereditary explanation would all be able to be improved as more genomes are sequenced and clarified. Additionally, current prescient calculations, for example, SIFT and Polyphen are reliant upon current data set explanation. On the off chance that enormous quantities of human genomes are sequenced, investigation need not hotel to only anticipating the impact of a solitary position; one can basically question that situation in the large numbers of individuals that are sequenced and induce the pernicious impact.

Other than the more unsurprising changes that will follow normally from more information, there are additionally openings for bigger perspective changes in bioinformatics instruments. In the first place, arising devices might have the option to examine tests not as a homogenous entire, but rather in manners that consider growth heterogeneity with varying populaces of cells. Besides, single-cell and single-atom strategies are developing. It is presently more liked that the cancers comprises of populaces of cells, and that having the option to decide the amount and personality of these cells won't just assist with understanding growth populace elements, yet may likewise illuminate therapy and anticipation.

Second, so far somewhat couple of instruments have incorporated other high throughput modalities, for example, proteomics into genomic translation. To comprehend whether the transformation has natural importance, it is basic to know whether a quality is communicated on a record or protein level. As more multidimensional information is created through tasks like ENCODE, TCGA, or 1000 Genomes, and high-throughput test profiling becomes simpler on a genomic, transcriptomic, and proteomic level, techniques that can fuse this information will add capacity to the investigation.

# Journal of Computational Biology

Third, in extra to multidimensional information, there are additionally openings for frameworks science techniques to be consolidated to programming bundles. Protein-protein collaboration datasets keep on developing as the human interactome is planned, and information on these sub-atomic pathways can and ought to be incorporated into genomics investigation. Understanding qualities as secluded develops as well as a component of a more prominent framework would better model the organic interaction.

Fourth, as increasingly more datasets are accessible and sequencing becomes less expensive, genomics examination need presently don't be founded on a solitary genome, a correlation between a separated pair of malignancy genome tests, or bigger, yet confined, families. Current instruments break down single examples all at once and contrast what is found and data sets. All things being equal, apparatuses that can investigate huge quantities of genomes simultaneously to sizes like genome-wide affiliation studies will end up being incredible.

Without a doubt, the datasets utilized in genomics examination will keep on filling top to bottom per individual and in the quantity of tests. Bioinformatics, like never before previously, will be the significant stage in sorting out the information flood. The gradual advancement managed by this flood will be basic and important, however specialists can likewise anticipate the yet-obscure outlook changes that loom into the great beyond.

## References

1. E. R. Mardis, "Next-generation DNA sequencing methods," *Annual Review of Genomics and Human Genetics*, vol. 9, pp. 387–402, 2008.View at: Publisher Site | Google Scholar

2. H. Li, B. Handsaker, A. Wysoker et al., "The Sequence Alignment/Map format and SAMtools," *Bioinformatics*, vol. 25, no. 16, pp. 2078–2079, 2009.View at: Publisher Site | Google Scholar

3. P. Danecek, A. Auton, G. Abecasis et al., "The variant call format and VCFtools," *Bioinformatics*, vol. 27, no. 15, Article ID btr330, pp. 2156–2158, 2011.View at: Publisher Site | Google Scholar

4. A. G. Day-Williams and E. Zeggini, "The effect of next-generation sequencing technology on complex trait research," *European Journal of Clinical Investigation*, vol. 41, no. 5, pp. 561–567, 2011.View at: Publisher Site | Google Scholar

5. M. Ruffalo, T. LaFramboise, and M. Koyuturk, "Comparative analysis of algorithms for next-generation sequencing read alignment," *Bioinformatics*, vol. 27, no. 20, pp. 2790–2796, 2011.View at: Google Scholar

6. N. Homer and S. F. Nelson, "Improved variant discovery through local re-alignment of short-read next-generation sequencing data using SRMA," *Genome Biology*, vol. 11, no. 10, article R99, 2010.View at: Publisher Site | Google Scholar

7. B. Langmead, C. Trapnell, M. Pop, and S. L. Salzberg, "Ultrafast and memory-efficient alignment of short DNA sequences to the human genome," *Genome Biology*, vol. 10, no. 3, article R25, 2009.View at: Publisher Site | Google Scholar

8. B. Langmead and S. L. Salzberg, "Fast gapped-read alignment with Bowtie 2," *Nature Methods*, vol. 9, no. 4, pp. 357–359, 2012.View at: Google Scholar

9. H. Li and R. Durbin, "Fast and accurate short read alignment with Burrows-Wheeler transform," *Bioinformatics*, vol. 25, no. 14, pp. 1754–1760, 2009.View at: Publisher Site | Google Scholar

10. H. Li, J. Ruan, and R. Durbin, "Mapping short DNA sequencing reads and calling variants using mapping quality scores," *Genome Research*, vol. 18, no. 11, pp. 1851–1858, 2008.View at: Publisher Site | Google Scholar

11. H. Li and R. Durbin, "Fast and accurate long-read alignment with Burrows-Wheeler transform," *Bioinformatics*, vol. 26, no. 5, Article ID btp698, pp. 589–595, 2010.View at: Publisher Site | Google Scholar

12. C. Alkan, S. Sajjadian, and E. E. Eichler, "Limitations of next-generation genome sequence assembly," *Nature Methods*, vol. 8, no. 1, pp. 61–65, 2011.View at: Publisher Site | Google Scholar

13. F. Hach, F. Hormozdiari, C. Alkan et al., "MrsFAST: a cache-oblivious algorithm for short-read mapping," *Nature Methods*, vol. 7, no. 8, pp. 576–577, 2010.View at: Publisher Site | Google Scholar

14. I. D. Dinov, F. Torri, F. Macciardi et al., "Applications of the pipeline environment for visual informatics and genomics computations," *BMC Bioinformatics*, vol. 12, article 304, 2011.View at: Publisher Site | Google Scholar

15. C. Alkan, J. M. Kidd, T. Marques-Bonet et al., "Personalized copy number and segmental duplication maps using next-generation sequencing," *Nature Genetics*, vol. 41, no. 10, pp. 1061–1067, 2009.View at: Publisher Site | Google Scholar

16. S. M. Rumble, P. Lacroute, A. V. Dalca, M. Fiume, A. Sidow, and M. Brudno, "SHRiMP: accurate mapping of short color-space reads," *PLoS Computational Biology*, vol. 5, no. 5, Article ID e1000386, 2009.View at: Publisher Site | Google Scholar

17. M. David, M. Dzamba, D. Lister, L. Ilie, and M. Brudno, "SHRiMP2: sensitive yet practical short read mapping," *Bioinformatics*, vol. 27, no. 7, Article ID btr046, pp. 1011–1012, 2011.View at: Publisher Site | Google Scholar

18. R. Li, Y. Li, K. Kristiansen, and J. Wang, "SOAP: short oligonucleotide alignment program," *Bioinformatics*, vol. 24, no. 5, pp. 713–714, 2008.View at: Publisher Site | Google Scholar

19. R. Li, C. Yu, Y. Li et al., "SOAP2: an improved ultrafast tool for short read alignment," *Bioinformatics*, vol. 25, no. 15, pp. 1966–1967, 2009.View at: Publisher Site | Google Scholar

20. C. M. Liu, K. F. Wong, E. M. K. Wu et al., "SOAP3: ultra-fast GPU-based parallel alignment tool for short reads," *Bioinformatics*, vol. 28, no. 6, pp. 878–879, 2012.View at: Google Scholar

21. R. Nielsen, J. S. Paul, A. Albrechtsen, and Y. S. Song, "Genotype and SNP calling from next-generation sequencing data," *Nature Reviews Genetics*, vol. 12, no. 6, pp. 443–451, 2011.View at: Publisher Site | Google Scholar

22. A. McKenna, M. Hanna, E. Banks et al., "The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data," *Genome Research*, vol. 20, no. 9, pp. 1297–1303, 2010.View at: Publisher Site | Google Scholar

23. Cancer Genome Atlas Network, "Comprehensive molecular portraits of human breast tumours," *Nature*, vol. 490, no. 7418, pp. 61–70, 2012.View at: Google Scholar

24. D. L. Altshuler, R. M. Durbin, G. R. Abecasis et al., "A map of human genome variation from population-scale sequencing," *Nature*, vol. 467, no. 7319, pp. 1061–1073, 2010.View at: Publisher Site | Google Scholar

25. D. C. Koboldt, K. Chen, T. Wylie et al., "VarScan: variant detection in massively parallel sequencing of individual and pooled samples," *Bioinformatics*, vol. 25, no. 17, pp. 2283–2285, 2009.View at: Publisher Site | Google Scholar

26. D. C. Koboldt, Q. Zhang, D. E. Larson et al., "VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing," *Genome Research*, vol. 22, no. 3, pp. 568–576, 2012.View at: Google Scholar

27. D. Challis, J. Yu, U. S. Evani et al., "An integrative variant analysis suite for whole exome next-generation sequencing data," *BMC Bioinformatics*, vol. 13, article 8, 2012.View at: Google Scholar

28. J. Hillman-Jackson, D. Clements, D. Blankenberg, J. Taylor, and A. Nekrutenko, "Using Galaxy to perform large-scale interactive data analyses," in *Current Protocols in Bioinformatics*, chapter 10, unit 10.5, 2012.View at: Google Scholar

29. H. P. Ji, "Improving bioinformatic pipelines for exome variant calling," *Genome Medicine*, vol. 4, no. 1, article 7, 2012.View at: Google Scholar

30. R. R. Lemos, M. B. Souza, and J. R. Oliveira, "Exploring the implications of INDELs in neuropsychiatric genetics: challenges and perspectives," *Journal of Molecular Neuroscience*, vol. 47, no. 3, pp. 419–424, 2012.View at: Google Scholar

31. S. A. Lee, H. S. Mun, H. Kim et al., "Naturally occurring hepatitis B virus X deletions and insertions among Korean chronic patients," *Journal of Medical Virology*, vol. 83, no. 1, pp. 65–70, 2011.View at: Publisher Site | Google Scholar

32. U. Väli, M. Brandström, M. Johansson, and H. Ellegren, "Insertion-deletion polymorphisms (indels) as genetic markers in natural populations," *BMC Genetics*, vol. 9, article 8, 2008.View at: Publisher Site | Google Scholar

33. G. Lunter, "Probabilistic whole-genome alignments reveal high indel rates in the human and mouse genomes," *Bioinformatics*, vol. 23, no. 13, pp. i289–i296, 2007.View at: Publisher Site | Google Scholar

34. P. Krawitz, C. Rödelsperger, M. Jäger, L. Jostins, S. Bauer, and P. N. Robinson, "Microindel detection in short-read sequence data," *Bioinformatics*, vol. 26, no. 6, Article ID btq027, pp. 722–729, 2010.View at: Publisher Site | Google Scholar

35. K. Ye, M. H. Schulz, Q. Long, R. Apweiler, and Z. Ning, "Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads," *Bioinformatics*, vol. 25, no. 21, pp. 2865–2871, 2009.View at: Publisher Site | Google Scholar

36. D. G. MacArthur, S. Balasubramanian, A. Frankish et al., "A systematic survey of loss-of-function variants in human protein-coding genes," *Science*, vol. 335, no. 6070, pp. 823–828, 2012.View at: Google Scholar

37. C. A. Albers, G. Lunter, D. G. MacArthur, G. McVean, W. H. Ouwehand, and R. Durbin, "Dindel: accurate indel calls from short-read data," *Genome Research*, vol. 21, no. 6, pp. 961–973, 2011.View at: Publisher Site | Google Scholar

38. J. A. Neuman, O. Isakov, and N. Shomron, "Analysis of insertion-deletion from deep-sequencing data: software evaluation for optimal detection," *Briefings in Bioinformatics*. In press.View at: Publisher Site | Google Scholar

39. G. M. Cooper and J. Shendure, "Needles in stacks of needles: finding disease-causal variants in a wealth of genomic data," *Nature Reviews Genetics*, vol. 12, no. 9, pp. 628–640, 2011.View at: Google Scholar

40. E. S. Lander, "Initial impact of the sequencing of the human genome," *Nature*, vol. 470, no. 7333, pp. 187–197, 2011.View at: Publisher Site | Google Scholar

41. L. A. Hindorff, P. Sethupathy, H. A. Junkins et al., "Potential etiologic and functional implications of genome-wide association loci for human diseases and traits," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 106, no. 23, pp. 9362–9367, 2009.View at: Publisher Site | Google Scholar

42. P. Kumar, S. Henikoff, and P. C. Ng, "Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm," *Nature Protocols*, vol. 4, no. 7, pp. 1073–1081, 2009.View at: Publisher Site | Google Scholar

43. I. A. Adzhubei, S. Schmidt, L. Peshkin et al., "A method and server for predicting damaging missense mutations," *Nature Methods*, vol. 7, no. 4, pp. 248–249, 2010.View at: Publisher Site | Google Scholar

44. P. S. Nair and M. Vihinen, "VariBench: A benchmark database for variations," *Human Mutation*. In press.View at: Google Scholar

45. P. Cingolani, A. Platts, L. Wang le et al., "A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w1118, iso-2, iso-3," *Fly*, vol. 6, no. 2, pp. 80–92, 2012.View at: Google Scholar

46. G. De Baets, J. Van Durme, J. Reumers et al., "SNPeffect 4.0: on-line prediction of molecular and structural effects of protein-coding variants," *Nucleic Acids Research*, vol. 40, pp. D935–D939, 2012.View at: Google Scholar

47. K. Wang, M. Li, and H. Hakonarson, "ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data," *Nucleic Acids Research*, vol. 38, no. 16, Article ID gkq603, p. e164, 2010.View at: Publisher Site | Google Scholar

48. M. Yandell, C. D. Huff, H. Hu et al., "A probabilistic disease-gene finder for personal genomes," *Genome Research*, vol. 21, no. 9, pp. 1529–1542, 2011.View at: Google Scholar

49. L. Habegger, S. Balasubramanian, D. Z. Chen et al., "VAT: a computational framework to functionally annotate variants in personal genomes within a cloud-computing environment," *Bioinformatics*, vol. 28, no. 17, pp. 2267–2269, 2012.View at: Google Scholar

50. D. G. MacArthur, S. Balasubramanian, A. Frankish et al., "A systematic survey of loss-of-function variants in human protein-coding genes," *Science*, vol. 335, no. 6070, pp. 823–828, 2012.View at: Google Scholar

51. M. Sincan, D. R. Simeonov, D. Adams et al., "VAR-MD: a tool to analyze whole exome-genome variants in small human pedigrees with mendelian inheritance," *Human Mutation*, vol. 33, no. 4, pp. 593–598, 2012.View at: Google Scholar